



Spurring Adoption of DHTs with OpenHash, a Public DHT Service

B. Karp, S. Ratnasamy, S. Rhea, S. Shenker

IRP-TR-03-16

November 10th, 2003

Research at Intel

INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL® PRODUCTS. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER, AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF INTEL PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT. Intel products are not intended for use in medical, life saving, life sustaining applications.

Intel may make changes to specifications and product descriptions at any time, without notice.

Copyright © Intel Corporation 2003

* Other names and brands may be claimed as the property of others.

Spurring Adoption of DHTs with OpenHash, a Public DHT Service

Brad Karp
Intel Research Pittsburgh
Carnegie Mellon

Sylvia Ratnasamy
Intel Research Berkeley

Sean Rhea
UC Berkeley

Scott Shenker
ICSI
UC Berkeley

1 Introduction

The past three years have seen intense research into Distributed Hash Tables (DHTs): both into algorithms for building them, and into applications built atop them. These applications have spanned a strikingly wide range, including file systems [5, 7, 11], event notification [12], content distribution [4], e-mail delivery [13], indirection services [1, 17], web caches [8], and relational query processors [10]. While this set of applications is impressively diverse, the vast majority of application building is done by a small community of DHT researchers. If DHTs are to have a positive impact on the design of distributed applications *used by real users outside this research community*, we believe that the community of DHT-based application developers should be as broad as possible.

Why, then, has this community of developers remained narrow? First, keeping a research prototype of a DHT running continually requires effort, and experience with DHT code. Second, significant testbed resources are required to deploy and test DHT-based applications. A hacker can download the code for Chord, but she cannot run that code alone; recall that only a tiny fraction of would-be developers has access to a testbed infrastructure like PlanetLab [14]. Consequently, most application developers would turn to ad hoc application-specific solutions rather than attempt to use a DHT.

Our central tenet is that we, as a community, need to harness the ingenuity and talents of the vast majority of application developers who reside outside the rarified but perhaps sterile air of the DHT research community. To that end, we issue a call-to-arms to deploy an *open, publicly accessible* DHT service that would allow new developers to experiment with DHT-based applications *without* the burden of deploying and maintaining a DHT. We believe that there are many simple applications that, individually, might not warrant the effort required to deploy a DHT but that would be trivial to build over a DHT service *were one available*. Many-to-many instant messaging and photo publication, where a user may share photos under a long-lived name even if the photos are served from a home machine with a dynamic IP address are but two of the many such applications. We term these *lite applications*, because they

make only simple and often fleeting use of a DHT.¹

To spur the development of these lite applications we propose *OpenHash*, an open, publicly accessible DHT service that runs on a set of infrastructure hosts and allows any application to execute `put()` and `get()` operations. Its presence would hopefully enable and encourage the development of a wide range of interesting and unexpected DHT applications.²

However, we already know from our limited experience with DHT-based applications that some require application-specific processing and thus can't be limited to using the generic `put()/get()` interface. The technical contribution of this paper is an examination of what one can do to extend the functionality of a DHT service so that such application-specific requirements can be reasonably accommodated. That is, we seek to share a common DHT routing platform while allowing application-specific functionality to reside only on certain nodes. To meet this hybrid challenge of shared routing but application-specific processing we propose ReDiR, a distributed rendezvous scheme that removes the need in today's DHT systems to place application-specific code together with DHT code. The ReDiR algorithm itself requires only a `put()/get()` interface from a DHT, but allows a surprisingly wide range of applications to work over a DHT service.

2 DHT as Library vs. as Service

Implementing a DHT as a service that exports a narrow `put()/get()` interface is a departure from how present-day DHT-based applications are built. To illuminate the consequences of this design decision, in this section we consider in turn the properties of the current approach and the service approach.

Many existing DHT applications (e.g. [4, 10, 11]) are built using a "bundled" model, where the application is able to read the local DHT state and receive upcalls from the DHT (as in [6]), either by being linked into the same process as the DHT code or through local RPC

¹Our claim is not that a DHT service is the *only* or *best* way to build these applications, but rather that a DHT service is a common building block that would be useful for a wide range of such applications, and would be far preferable to building one-off, special-purpose rendezvous or indirection mechanisms for each application (e.g., dynamic DNS for photo publication).

²We are often reminded that the most successful peer-to-peer application was developed by a 19-year-old.

calls. For brevity, we will say these applications use the DHT as a *library* in either case. To support upcalls, the library model requires that code for the same set of applications be available at all DHT hosts; in practice, this limitation prevents the sharing of a single DHT deployment. Instead, different applications only re-use the DHT *code*, amortizing the development of the DHT functionality over the set of applications that use it. A side effect of this is that every such application imposes the maintenance traffic associated with running a DHT on its underlying infrastructure.

The great strength of the library model is the flexibility it affords applications in functionality. By bundling arbitrary application code with DHT code, the model supports any operation at any overlay host on any data held at that host. However, this flexibility comes at the expense of synergy in deployment, and thus at the expense of ease of use of the DHT. We believe that these weaknesses in the library model are the primary reasons for the narrowness of the community of developers of DHT-based applications today. Most would-be DHT-using developers look elsewhere for rendezvous and indirection solutions, because the effort and resources required to deploy a DHT are often greater than those required to hack up a more ad hoc solution.

In contrast, a DHT that provides a more limited interface, specifically only `put()` and `get()`, does not need to be bundled with application code and can thus be deployed as a common, openly accessible service. There are two principal weaknesses of this service model. First, a publicly accessible service is subject to attack.³ We include a brief discussion in Section 4 but defer a full exploration of the problem to later work. The second problem, which is the focus of this paper, is that the service model is far less flexible in the functionality (and thus the applications) it supports. Intuitively the `put()/get()` interface appears far more restrictive than the “any code, any node” approach of the library model. We offer a detailed discussion of what the service model can and cannot support in the following section. The library model and service model are extremes on a continuum between maximal flexibility without an easily adopted infrastructure, and providing such an infrastructure with reduced flexibility. In the next section, we explore the space between these extremes and the extent to which application-specific functionality can live atop a DHT service.

3 Application-Specific Functionality in OpenHash

Existing, library-based DHT applications [11, 10, 17, 7] co-locate application-specific functionality with the DHT. We distinguish between application functional-

ity invoked at the endpoint (*i.e.*, destination) of a DHT route versus functionality invoked at *every* hop along a route and show that, surprisingly, the former is quite easily achievable over a service while the latter is not.

We begin this section by illustrating the distinction between endpoint and per-hop functionality in Section 3.1; we describe one solution to achieving endpoint functionality over a service in Section 3.2 and end with a discussion of per-hop functions in Section 3.3.

3.1 Endpoint vs. Per-Hop Operators

We use illustrative examples drawn from existing applications, namely PIER and *i3*, to expose the difference between endpoint and per-hop functionality.

The first system we consider is *i3*. In *i3*, packets are *forwarded* to identifiers called triggers. To send a packet to x , a sender routes the packet to the DHT host responsible for storing trigger entries for x . That DHT host then extracts the value I corresponding to x and *forwards* the packet to I . There are two aspects to *i3*’s support for packet forwarding within a DHT. First, each DHT host in *i3* must implement the forwarding operation. Second, that forwarding code must read the (key, value) pair stored locally at that DHT host.

Likewise, to execute a join operation in PIER [10], a DHT host iterates over all the (key, value) pairs in its local store, and rehashes them by a portion of their value fields. Here the code for a join operator resides at every DHT host and has full access to the host’s local store. The essence of the above behaviors is that the DHT routes an operation request to a key within the keyspace, and the end DHT host responsible for that key carries out an operation that typically accesses the key-value entries stored locally. We term the combination of these behaviors an *endpoint operator*.

A somewhat different example is PIER’s computation of aggregates along a tree rooted at a particular key. Nodes route messages toward the root and each DHT host along the path aggregates data before forwarding them. Multicast forwarding within a DHT (as in Scribe, Bayeux, or M-CAN) also uses per-hop processing to set up and maintain dissemination trees. We term such behavior a *per-hop operator*, as it requires application-specific operations be executed at each hop along the path to a key (as opposed to at the final node that holds the key).

3.2 Supporting Endpoint Operators in OpenHash

Note that OpenHash cannot by itself support the implementation of either endpoint or per-hop operators because, as a shared service, these operators do not reside at the hosts that constitute OpenHash’s DHT. However, we can allow application specific code to live *outside* of OpenHash and use OpenHash to direct requests to nodes that do support the required operators. This ap-

³Note, however, that almost any deployed DHT-based system is subject to attack, whether it uses a DHT library or service.

proach allows developers to deploy application-specific code at will while still using OpenHash’s common routing infrastructure. Note that as with any DHT, an application will now need to split its keyspace over the application nodes. In other words, when routing a request to key k for application A , the routing system must find the host that runs application A and is responsible for key k in A ’s keyspace. Conventional DHTs consistently hash a set of keys over the hosts that run the DHT itself. We need to solve a different problem: to consistently hash a set of keys over a set of application hosts that do not run the DHT itself. This functionality is effectively the same as that of the `lookup()` interface first proposed in the Chord paper and adopted by the authors in [6] as the KBR or Key-Based Routing interface, with one important difference: our `lookup()` maps application keys to arbitrary application hosts, rather than to the hosts that run the DHT.

In this section, we describe ReDiR (*Recursive Distributed Rendezvous*), a mechanism by which OpenHash can be used to achieve such application-specific `lookup()`s.⁴ ReDiR requires hosts that support endpoint operators for an application to register this with OpenHash. Clients can then use OpenHash to route requests for a particular key and a particular application to the application host responsible for that key.

ReDiR: ReDiR hierarchically decomposes the OpenHash keyspace into nested partitions. An example of such a binary decomposition is shown in Figure 1(a)—level 0 in the decomposition corresponds to the entire keyspace $[0 \dots K]$ while at level l the keyspace is decomposed into 2^l partitions, each of length $K/(2^l)$.⁵ Thus, every point in the keyspace has a corresponding set of enclosing partitions, one at every level in the decomposition tree. To rendezvous on a name (say ABC) in a localized⁶ manner a node X first computes $H(X)$, which determines its set of enclosing partitions. It then hashes the name ABC within the scope of its smallest enclosing partition, and continues to do so within successively larger enclosing partitions until it finds the desired rendezvous information. Figure 1(a) shows the enclosing partitions and multi-level rendezvous points for node X .

ReDiR supports operators in ABC as follows:

- To join an application ABC , a node X performs a `get(ABC)` within its successively larger enclosing partitions until it finds an entry of the form

⁴Note that ReDiR solves a very different problem from that in [3]. There the authors propose a single DHT used to bootstrap multiple application-specific DHTs; we’re proposing a single DHT shared by multiple applications.

⁵Note that the base and depth of decomposition need not be the same across applications, and could in fact be encoded into the rendezvous name.

⁶Locality here refers to locality in keyspace.

$(H(ABC), Y)$ such that Y is a predecessor of X in the keyspace (i.e., $H(Y) < H(X)$ with wrap-around). Let k be the level at which this predecessor entry is located.

- Starting from the deepest level, X now inserts a key-value entry of the form $H(ABC, X)$ at each level in the decomposition until it reaches level k .
- To locate the successor for a key k among the nodes in application ABC , a node computes $H(k)$ and performs `get(ABC)` in successively larger enclosing partitions corresponding to $H(k)$ until it finds an entry of the form $(H(ABC), Z)$ such that $k < H(Z)$ at which point it has located the required successor.

Figure 1(b) shows an example of the above. We point out that the entire ReDiR mechanism builds entirely over the normal DHT `put()/get()` interface. To OpenHash, the ReDiR-related (key, value) pairs appear as any others.

Even with ReDiR, OpenHash leaves developers with the burden of deploying application-specific operators. We imagine that over time OpenHash will grow to incorporate some of these more specialized operators, but we don’t yet know what this subset should be. Moreover, expecting every node in the OpenHash infrastructure always to run exactly the same set of operators is unrealistically utopian. ReDiR enables a single routing layer to be shared by all services whether partially deployed or not. Without ReDiR, we are stuck with either utopia or bust.

3.3 Supporting Per-hop Operators in OpenHash

Although a DHT service with ReDiR, performs the route-to-key function for an application, there is no obvious way to support per-hop operators in OpenHash. However, in considering the various forms of per-hop operators described in the literature, we discovered that in most cases one could achieve similar functionality outside of the DHT service, essentially by converting per-hop operations into “scoped” endpoint operations (which ReDiR supports). While we do not claim to know whether all per-hop operators can be implemented using a service, we find this approach promising, and will explore it in greater detail in future work. In this section, we very briefly describe how three sample operations that typically use per-hop operations—multicast, aggregation and server selection (DOLR)—can be built over a DHT service.

Multicast: The following is a brief sketch of one possible solution to multicast in which OpenHash (using ReDiR) provides the rendezvous mechanism, while end nodes implement forwarding. To join group G , a node A inserts (G, A) within successively larger partitions enclosing $H(A)$ until it hits a partition in which there is already an entry for G . Let d be the maximum depth of

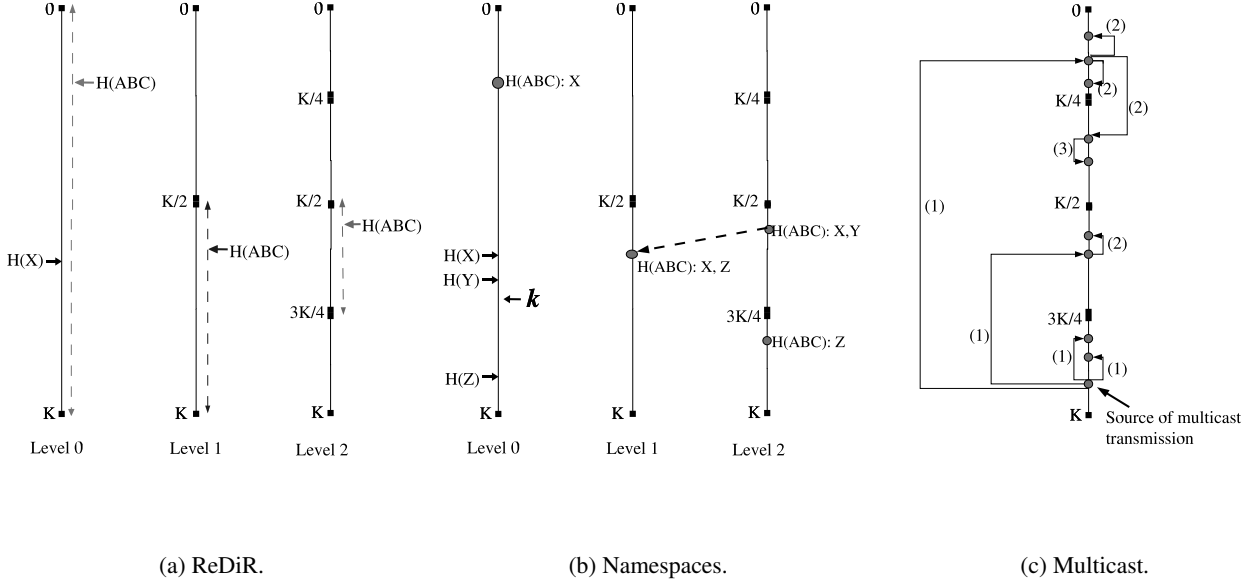


Figure 1: 1(a) ReDiR hierarchy and rendezvous points for a node X belonging to application ABC . Dashed lines denote enclosing partitions for X at each level. For clarity, the keyspace at each level is redrawn; in practice there is only one keyspace. 1(b) Using ReDiR to find the correct successor node for application ABC 's key k : Node X joined first, then Y and Z . Each rendezvous point lists the nodes registered there. Node Z is key k 's correct successor. Dashed lines denote rendezvous nodes contacted to locate k 's successor. 1(c) Example of multicast forwarding: transmission links are labelled with their depth in forwarding hierarchy.

the ReDiR hierarchy. Then, to multicast to all of group G , node A does a `get(G)` within its level d partition and unicasts the message to all the group members returned by the `get()`. A also does a `get(G)` within each of its “sibling” partitions from level $d - 1$ to 0 and unicasts a message to any *one* node at each level. A node, say B , that receives a message from A assumes that it is responsible for forwarding it on within its half of the smallest partition in the decomposition that contains both A and B ; if that decomposition is at level d , B does nothing. Figure 1(c) shows an example of this forwarding. Note however, that unlike schemes like Scribe[4], the above solution need not result in trees optimized for low latency.

Aggregation: Aggregation along the path to a root R could be implemented similarly to multicast by having rendezvous nodes for R at level i aggregate messages before forwarding them on to level $i - 1$.

Server selection using DOLR: In systems such as OceanStore and PAST, a client's lookup returns the address of the node closest to the client that stores a copy of the requested object. These systems achieve this through a combination of proximity-sensitive DHT construction and by caching pointers along the path between a node storing a copy and the root node for that object's identifier. This mechanism (often called

DOLR) serves two purposes: (1) server selection based on network latency and (2) fate sharing in the sense that if a closeby (e.g., within the client's organization) server is available, the lookup will succeed even if a large fraction of the DHT is unavailable to the client. Both of the above can be achieved without explicitly embedding the supporting functionality into the DHT. For example, [1] and [13] use an org-store and local rings to achieve fate sharing, and [5, 17] use network coordinates to find close copies. In fact, such approaches frequently give applications more control over selection criteria (e.g., server bandwidth or load could easily be used in place of or in addition to latency). Such flexibility is much harder to achieve using the DOLR approach.

4 Architecture Details

In this section we discuss architectural issues in the design of OpenHash. We begin with the service model, then discuss issues of resource contention.

4.1 Basic Service Model

Figure 2 presents an overview of the OpenHash architecture. Each PlanetLab host runs the OpenHash code and maintains a local store of the (key, value) pairs for which it is responsible. These local stores are accessible to OpenHash clients only through the `put()`/`get()` interface. Lite applications use only

- [16] T. Roscoe and S. Hand. Palimpsest: Soft-capacity storage for planetary-scale services. In *HotOS IX*. May 2003.
- [17] I. Stoica, D. Adkins, et al. Internet Indirection Infrastructure. In *Proceedings of the ACM SIGCOMM 2002*. Pittsburgh, PA, USA, August 2002.